**NTNU – Trondheim**
Norwegian University of
Science and Technology

# Is (Deep) Neural Networks Trustworthy?

Hao Wang, Big Data Lab, IIR, NTNU, Norway

# On my background

- Ph.D. and B.Eng. in Computer Engineering
- 6 years in Canada
  - 4 years in Nova Scotia (St. Francis Xavier Univ.)
    - Workflow modeling and HPC
  - 2 years in IBM Canada/McMaster Univ.
    - Software Certification, Safety Analysis
- Moved to Norway in 2014
  - Big Data Lab
    - Applied R&D and Education on Big Data (Visual Analytics) and IIoT
    - Applications from Ocean (Maritime and Aquaculture), Sustainability (Water and Air pollution), Health (Cancer), FinTech

# NTNU Ålesund

- Closely integrated with local industry
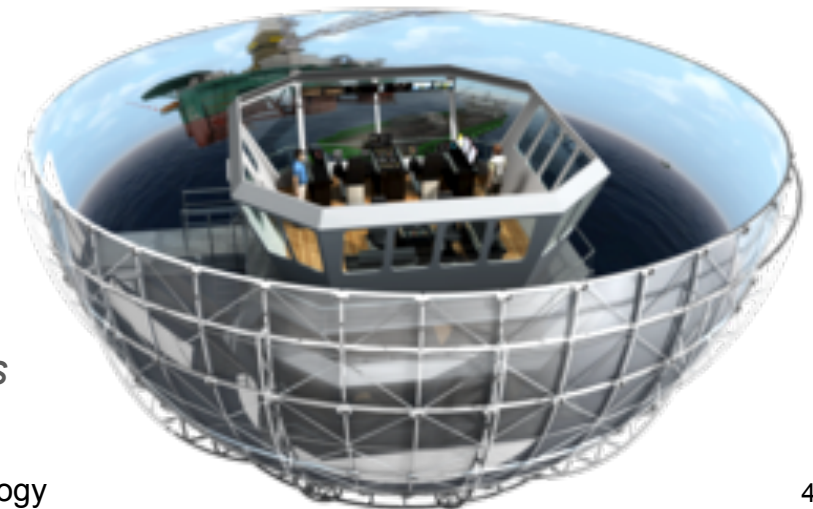- Strong interdisciplinary collaborative atmosphere

# Sunnmøre Maritime Cluster

**A world-leading maritime cluster**

- 13 design companies

  14 ship yards

  20 ship-owner companies

  169 equipment suppliers

  22 500 employees

  40% of the world's modern fleet



*Most of Norway's strength in the field of advanced marine operations is concentrated within an hour's drive from Ålesund. The region is home for an impressive constellation of over 200 leading maritime companies and training, research and finance institutions that form one of the most complete maritime clusters in the world.*
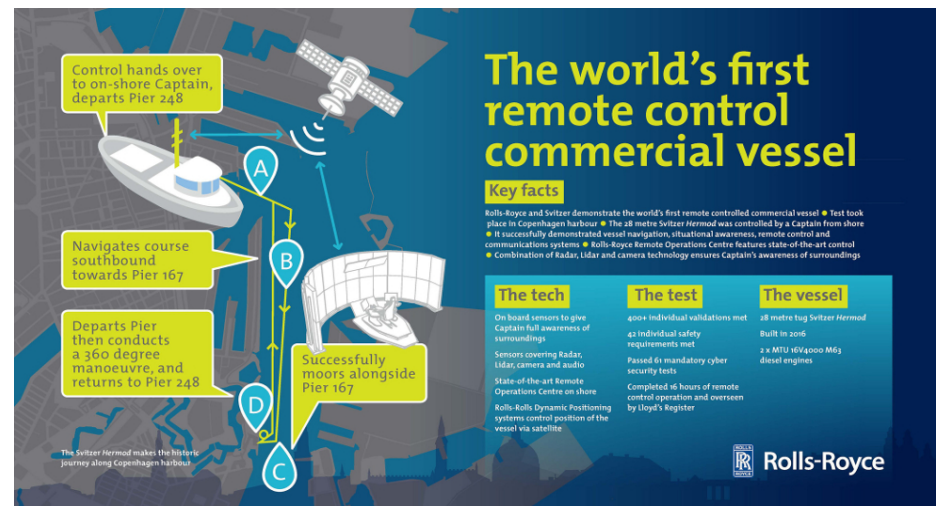
# Data Analytics for Maritime Operations

- Sensor data obtained through *HEalth MOnitoring System* (HEMOS) by Rolls Royce Marine AS
  - Time series + GPS
  - High frequency
    - E.g., vibration, torque
  - Low frequency
    - E.g., heading, speed, and GPS

# Security for Cyber-enabled Vessels

- A project with Profs. Sokratis K. Katsikas, Slobodan Petrovic, and Edmund Førland Brekke

- Goals
  - To analyze the cyber security risks of the cyber-enabled remotely-operated ship
  - To propose effective and efficient risk management strategies

June 2017: Rolls-Royce and Svitzer demonstrate world's first remotely operated commercial vessel in Copenhagen habour

# Robustness of Deep Neural Networks

- NN are being integrated in Safety Critical Systems
- Szegedy et al. [1] observed that state- of-the-art NN are highly vulnerable to *adversarial perturbations*
  - Given a correctly-classified input x, it is possible to find a new input x′ that is very similar to x but is as- signed a different label.
  - For instance, in image-recognition networks it is possible to add a small amount of noise (undetectable by the human eye) to an image and change how it is classified by the network [2].

[1] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus. Intriguing properties of neural networks, 2013. Technical Report.
[2] Divya Gopinath, Guy Katz, Corina S. Pasareanu, and Clark Barrett. DeepSafe: A Data-driven Approach for Checking Adversarial Robustness in Neural Networks. 2017

# Hypotheses on this Weakness

- 1. high complexity and non-linearity of neural networks can assign random labels in areas of the space which are under-explored [1].
  - Has been refuted
    - Unable justify the transferability of adversarial samples from one model to another
    - Linear models also suffer from this phenomenon

[1] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus. Intriguing properties of neural networks, 2013. Technical Report.

# Hypotheses on this Weakness

- 2. Goodfellow et al. [3] proposed a *linearity hypothesis* instead:

  – deep neural networks are highly non-linear with respect to their parameters, but mostly linear with respect to their inputs, and adversarial examples are easy to encounter when exploring a direction orthogonal to a decision boundary.

  – Another conjecture to explain the existence of adversarial examples is the cumulation of errors while propagating the perturbations from layer to layer.

    - A small carefully crafted perturbation in the input layer may result in a much greater difference in the output layer, effect that is only magnified in high dimensional spaces, causing the activation of the wrong units in the upper layers

[3] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. 2014.

NTNU    Norwegian University of Science and Technology

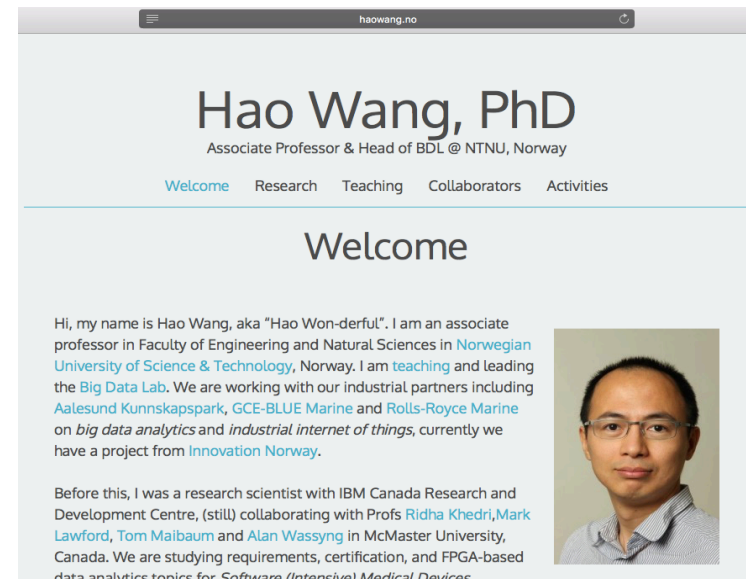# Approaches

- Casting to solve an optimization problem

$$\min_{\Delta x} F(x + \Delta x) \neq F(x), \quad \text{s.t.} \quad x' \in [0, 1]^n.$$

- Verification approach
  - E.g., Deep Learning Verification (DLV) is an approach that defines a region of safety around a known input and applies SMT solving for checking robustness

[4] X. Huang, M. Kwiatkowska, S. Wang, and M. Wu. Safety verification of deep neural networks. In Proc. 29th Int. Conf. on Computer Aided Verification (CAV), pages 3–29, 2017

# Tusen takk!!

- Comments?
- Collaborative ideas?



- Hao Wang, hawa@ntnu.no, www.haowang.no
- Big Data Lab 🔊 http://blog.hials.no/bigdata/